

Generación de trayectorias robóticas mediante aprendizaje profundo por refuerzo

MAYO 2018

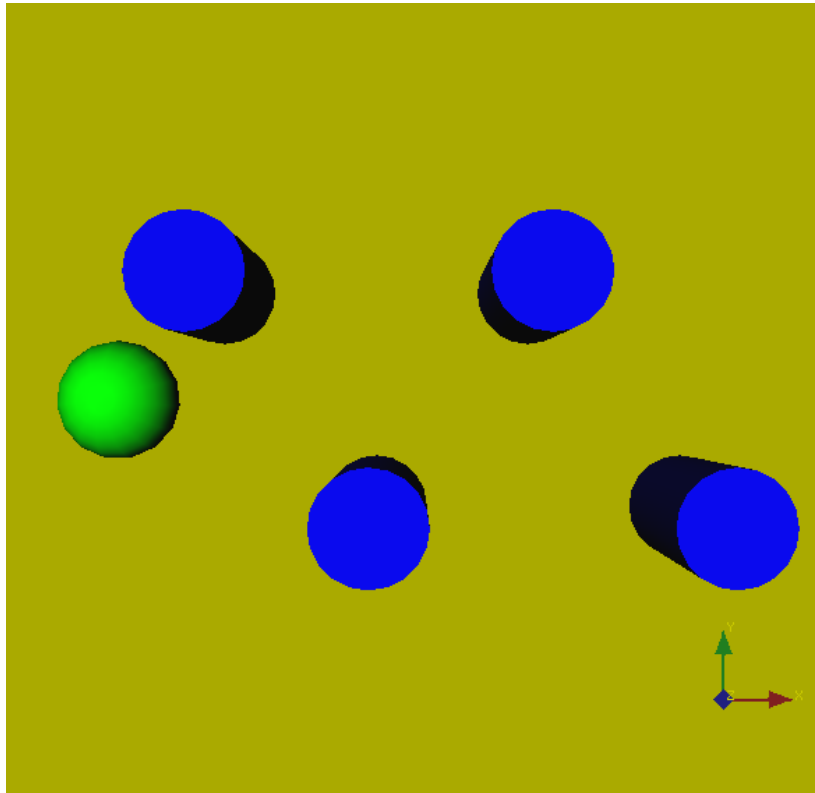
GUILLERMO URCERA MARTÍN

DIRECTOR: JAN ROSELL GRATACÒS





Problema objetivo



Elemento	Descripción
Dimensiones del área de trabajo	$19 \times 38\text{cm}$
Dimensiones de los obstáculos	Cilindros de 11,5cm de altura, 6,6cm de diámetro
Dimensiones del TCP	Esfera de 7cm de diámetro



Índice

Introducción

- Objetivos
- Motivación

Algoritmo

- Antecedentes
- Puntos clave
- Núcleo

Implementación

- Entornos
- Arquitectura
- Detalles
- Experimentos

Resultados

- Pruebas con modificaciones
- Comparativa con KPIECE

Conclusiones y trabajo futuro

Objetivos del trabajo

Implementar un sistema de planificación de movimientos utilizando el algoritmo DDPG

Comparar el rendimiento del sistema con otro planificador en simulación

Implementar el sistema en un robot real

Planificación de movimientos

Enfoque clásico:

- Métodos basados en la descomposición por celdas
- Métodos basados en campos de potenciales
- Métodos basados en mapas de ruta

Métodos basados en el muestreo

- Probabilistic Road Maps (PRM)
- Rapidly-exploring Random Trees (RRT)
- Kinodynamic Motion Planning by Interior-Exterior Cell Exploration (KPIECE)

Problemas de los planificadores convencionales

Alto coste computacional / temporal

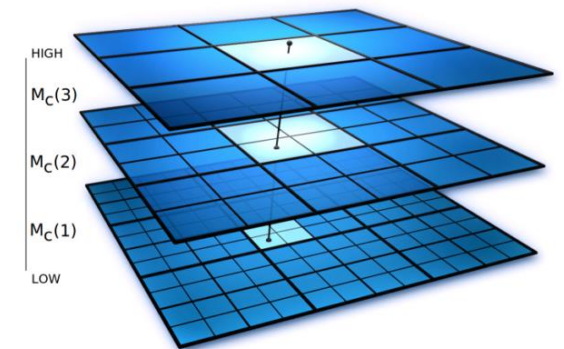
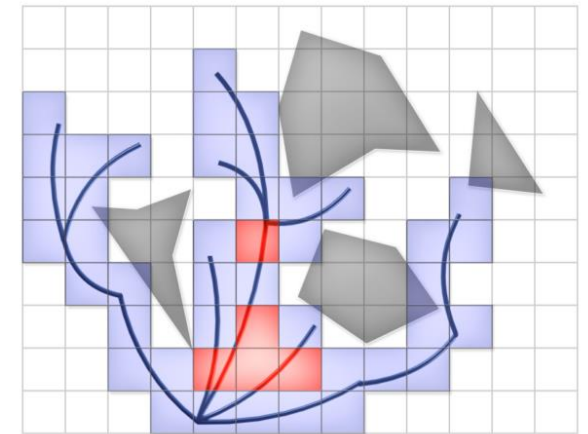
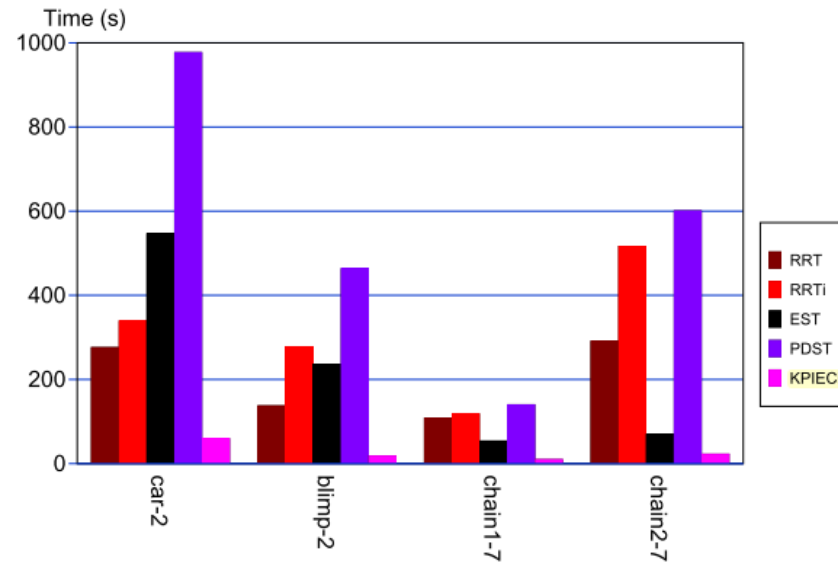
Curse of dimensionality

Precisión subordinada a una simulación física

Superficies del espacio de configuraciones complejas

Kinodynamic Motion Planning by Interior-Exterior Cell Exploration (KPIECE)

- Basado en el muestreo
- Tiempos de cálculo y uso de memoria reducidos
- Planificación con simulación física
- Exploración eficaz:
 - Proyecciones a espacios de menor dimensionalidad
 - Discretización a varios niveles



Aprendizaje automático

Las redes neuronales permiten escapar del *curse of dimensionality*

Las redes neuronales permiten generalizar

Mayor adaptabilidad al entorno

- Parámetros iniciales
- Cambios

Algoritmo

Introducción

- Objetivos
- Motivación

Algoritmo

- Antecedentes
- Puntos clave
- Núcleo

Implementación

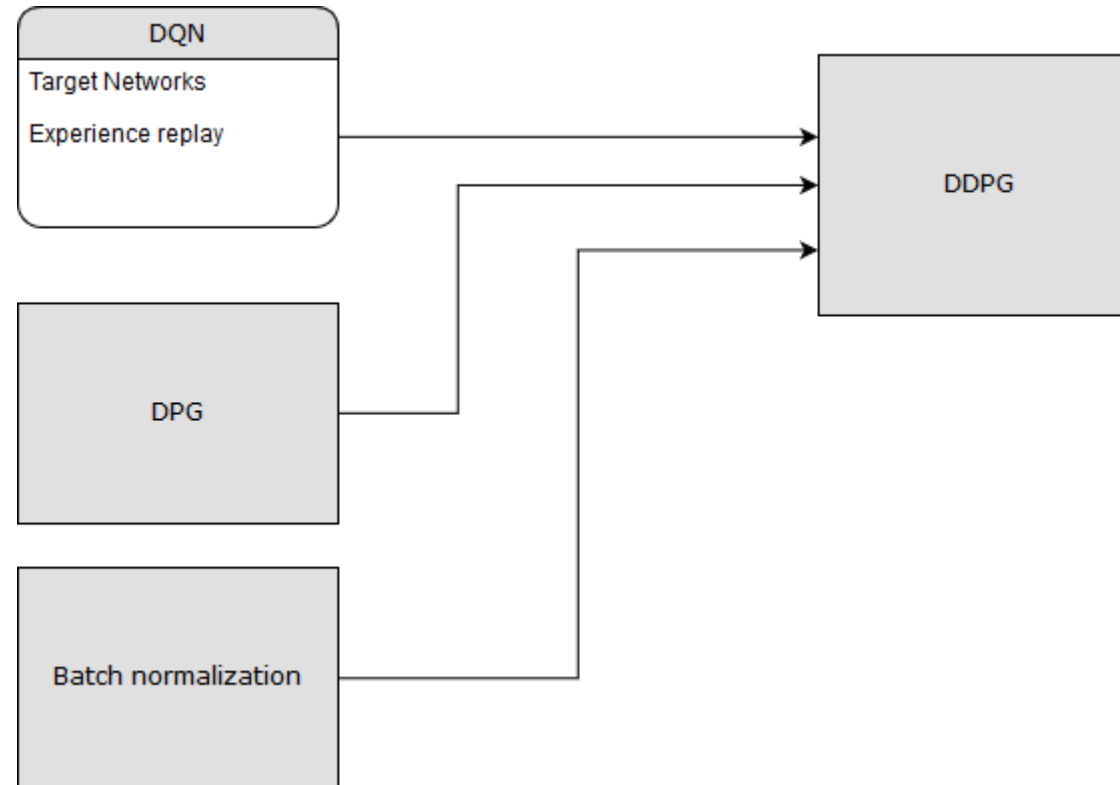
- Entornos
- Arquitectura
- Detalles
- Experimentos

Resultados

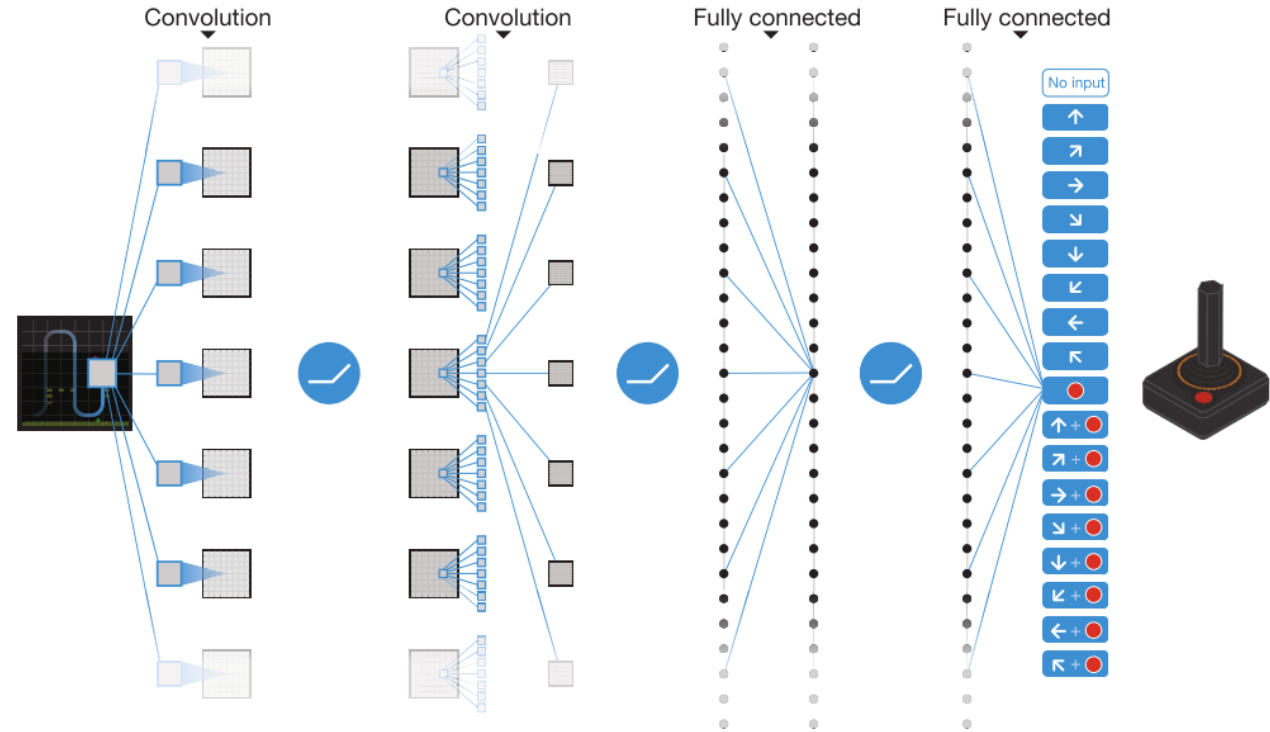
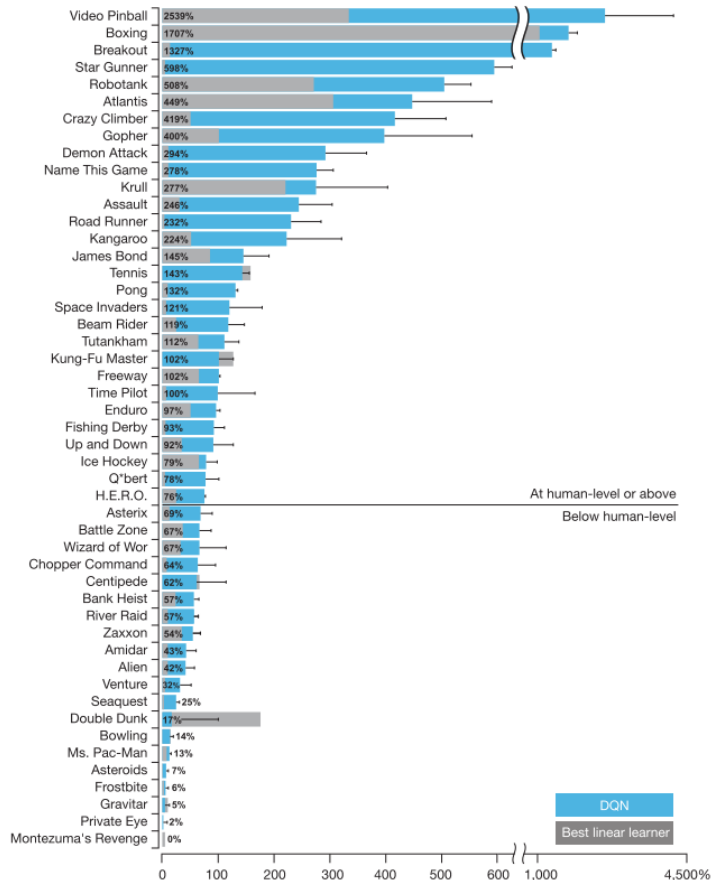
- Pruebas con modificaciones
- Comparativa con KPIECE

Conclusiones y trabajo futuro

Deep Deterministic Policy Gradient (DDPG)



Deep Q Network (DQN)



Deterministic Policy Gradient (DPG)

El gradiente de la política solo debe integrar sobre el espacio de estados

Rendimiento elevado en comparación con los gradientes de políticas estocásticas

$$\begin{aligned}\nabla_{\theta} J(\mu_{\theta}) &= \int_{\mathcal{S}} \rho^{\mu}(s) \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu}(s, a) \big|_{a=\mu_{\theta}(s)} ds \\ &= \mathbb{E}_{s \sim \rho^{\mu}} [\nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu}(s, a) \big|_{a=\mu_{\theta}(s)}]\end{aligned}$$

$$\begin{aligned}\nabla_{\theta} J(\pi_{\theta}) &= \int_{\mathcal{S}} \rho^{\pi}(s) \int_{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(a|s) Q^{\pi}(s, a) da ds \\ &= \mathbb{E}_{s \sim \rho^{\pi}, a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi}(s, a)]\end{aligned}$$

Puntos clave

Target networks

- Estabiliza el aprendizaje

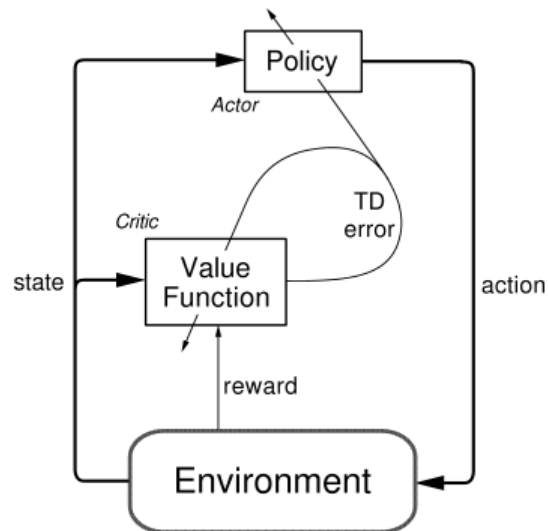
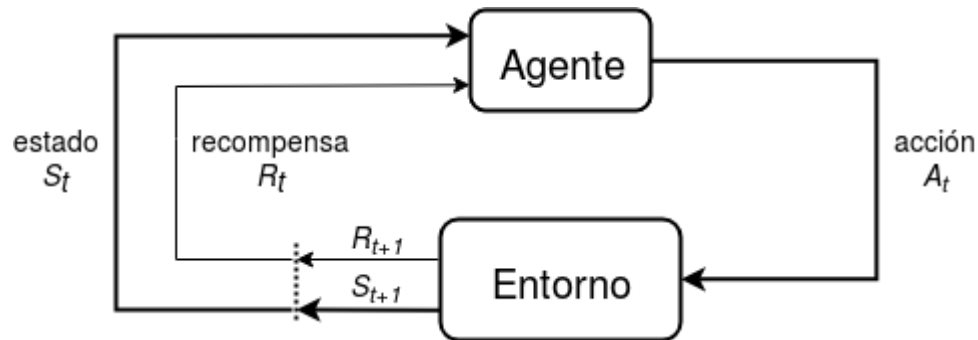
Experience replay

- Acelera el proceso de aprendizaje
- Aumenta la eficiencia de la exploración
- Estabiliza el aprendizaje

Batch normalization

- Reduce el *covariate shift*
- Acelera el aprendizaje

Deep Deterministic Policy Gradient (DDPG)



Algoritmo 1: Algoritmo DDPG[26]

```

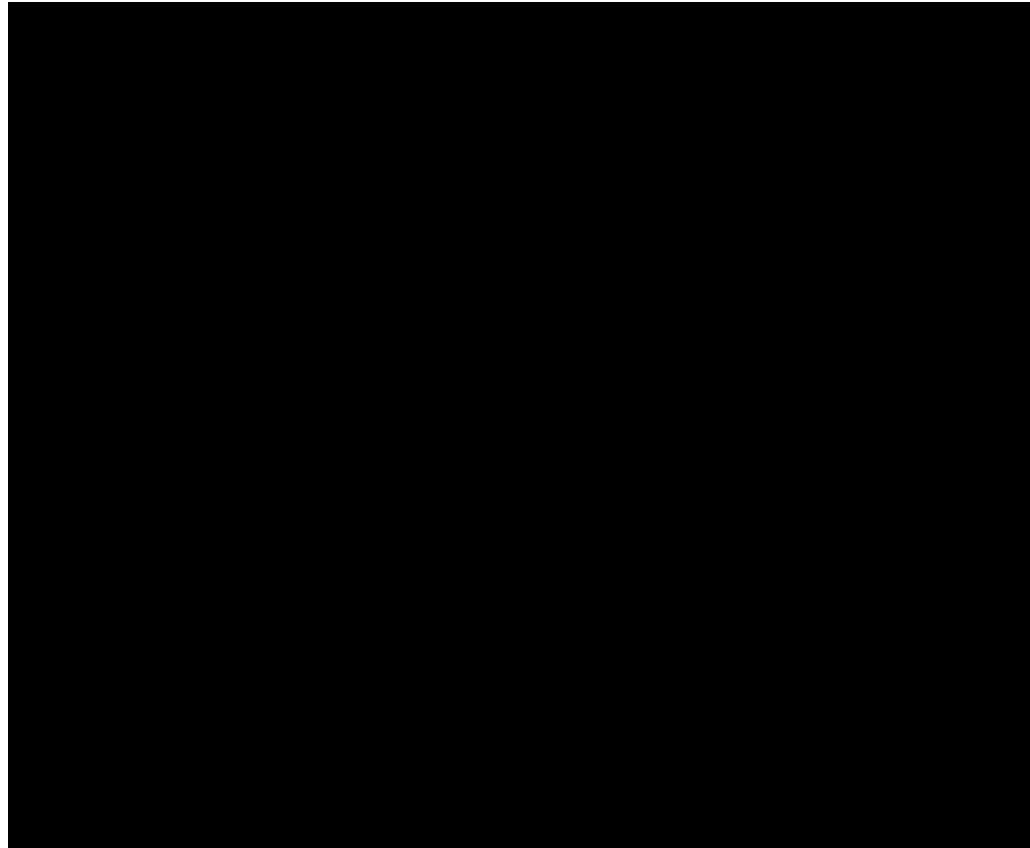
1 Inicializar aleatoriamente la red del crítico  $Q(s, a | \theta^\mu)$  y el actor
   $\mu(s | \theta^\mu)$  con pesos  $\theta^Q$  y  $\theta^\mu$ .
2 Inicializar las redes objetivo  $Q'$  y  $\mu'$  con pesos  $\theta^{Q'} \leftarrow \theta^Q, \theta^{\mu'} \leftarrow \theta^\mu$ 
3 Inicializar el buffer de interacciones  $R$ 
4 for episodio = 1,  $M$  do
5   Inicializar proceso aleatorio  $N$  para la exploración del espacio de
     acciones
6   Recibir estado observacional inicial  $s_1$ 
7   for  $t=1, T$  do
8     Seleccionar acción  $a_t = \mu(s_t | \theta^\mu) + N_t$  de acorde con la política
       actual y el ruido exploratorio
9     Ejecutar acción  $a_t$  y observar la recompensa  $r_t$  y el nuevo
       estado  $s_{t+1}$ 
10    Almacenar transición  $(s_t, a_t, r_t, s_{t+1})$  en  $R$ 
11    Muestrear un minibatch aleatorio de  $N$  transiciones
        $(s_i, a_i, r_i, s_{i+1})$  de  $R$ 
12    Sea  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1} | \theta^{\mu'}) | \theta^{Q'})$ 
13    Actualizar el crítico minimizando la pérdida:
       
$$L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i | \theta^Q))^2$$

14    Actualizar la política del actor utilizando el gradiente de la
       política muestreado:
       
$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a | \theta^Q) |_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s | \theta^\mu) |_{s_i}$$

15    Actualizar las redes objetivo:
16     $\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$ 
17     $\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$ 
18   end
19 end

```

Deep Deterministic Policy Gradient (DDPG)



Implementación

Introducción

- Objetivos
- Motivación

Algoritmo

- Antecedentes
- Puntos clave
- Núcleo

Implementación

- Entornos
- Arquitectura
- Detalles
- Experimentos

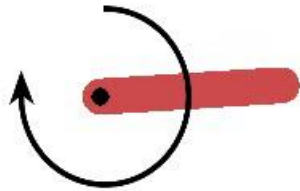
Resultados

- Pruebas con modificaciones
- Comparativa con KPIECE

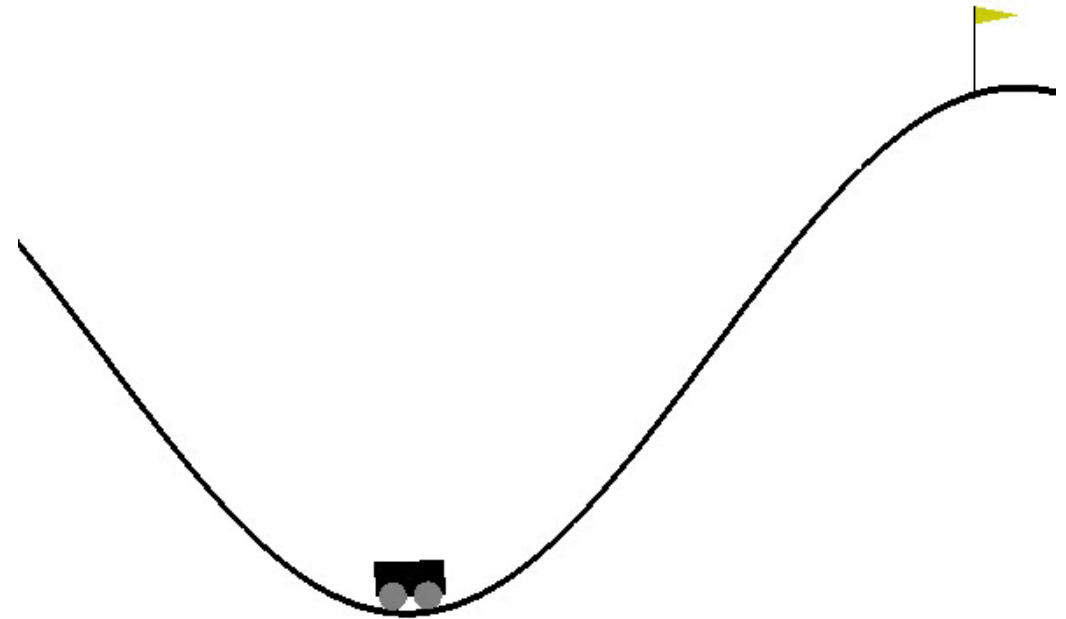
Conclusiones y trabajo futuro

Implementación en problemas sencillos

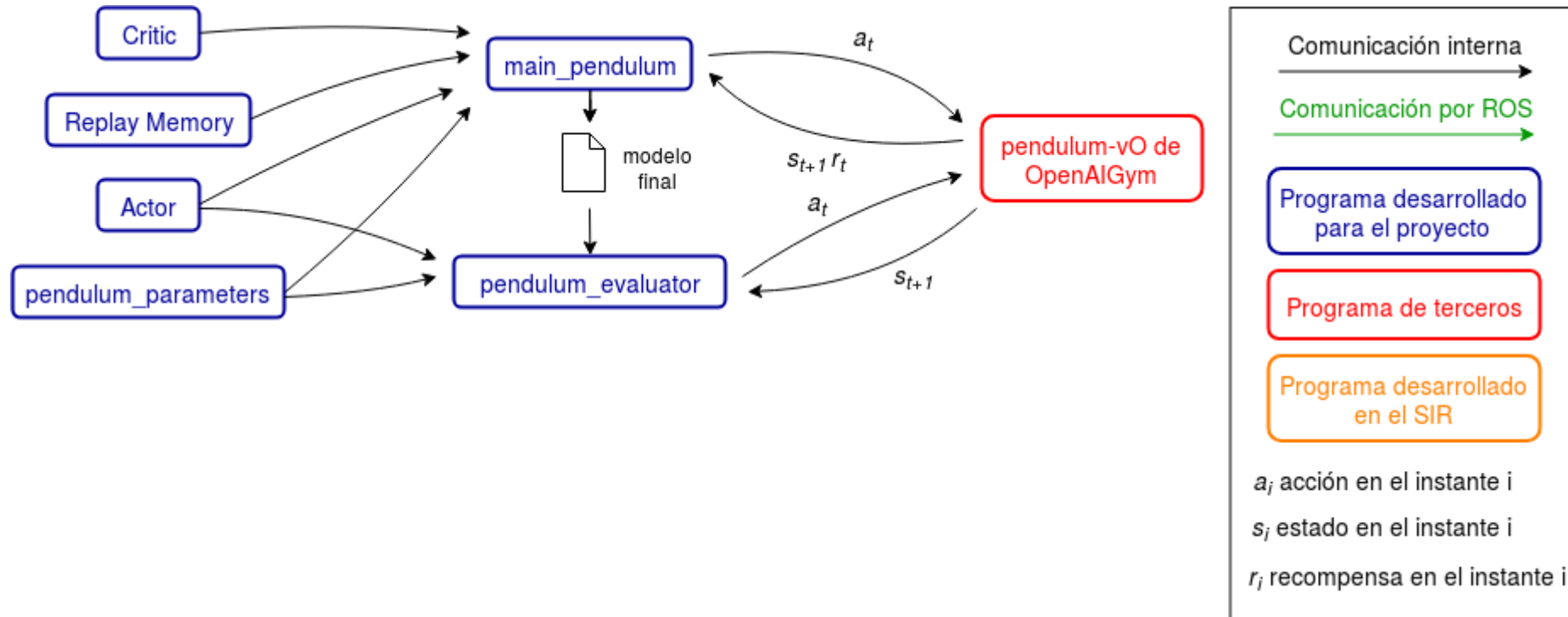
PENDULUM-V0



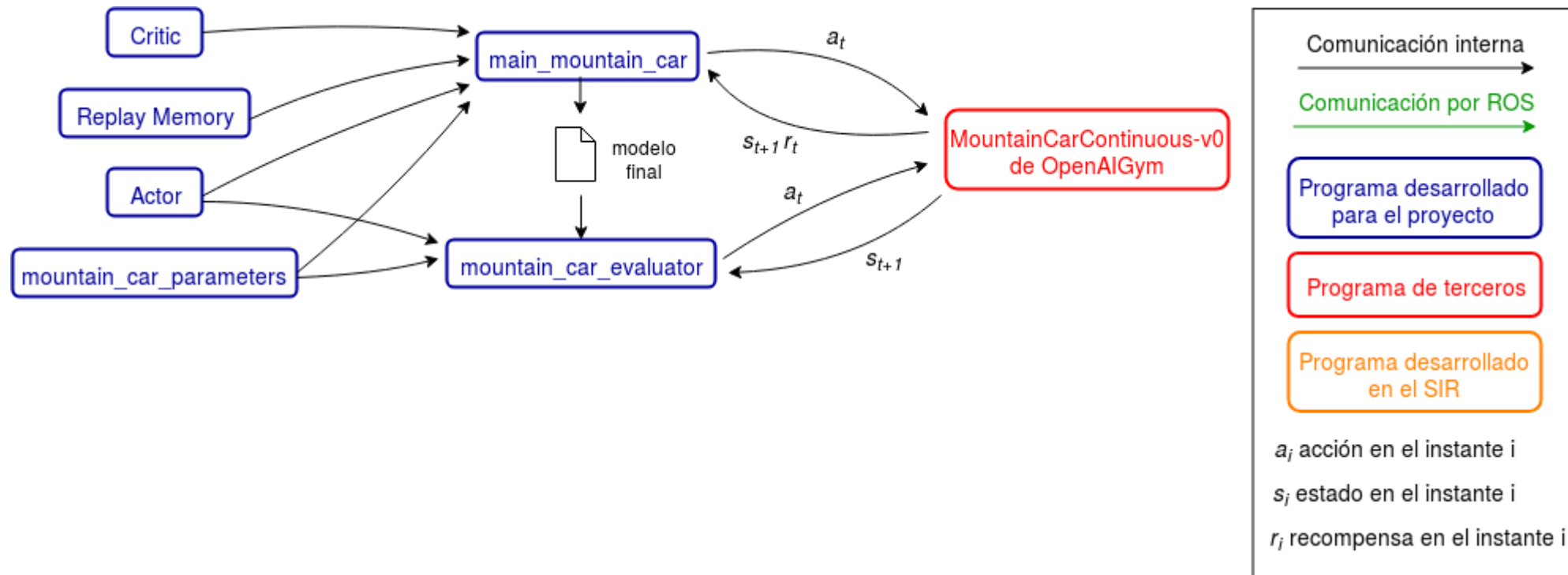
MOUNTAINCARCONTINUOUS-V0



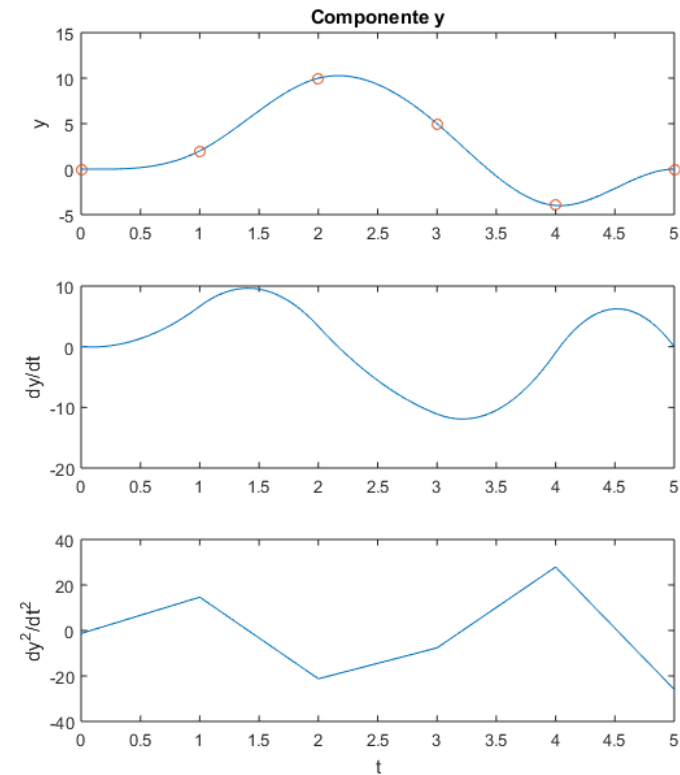
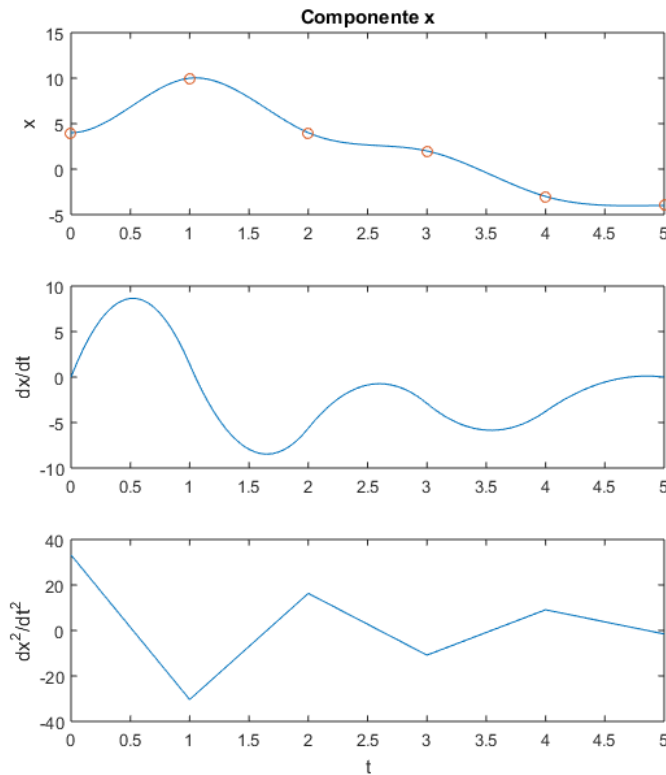
Pendulum



MountainCarContinuous

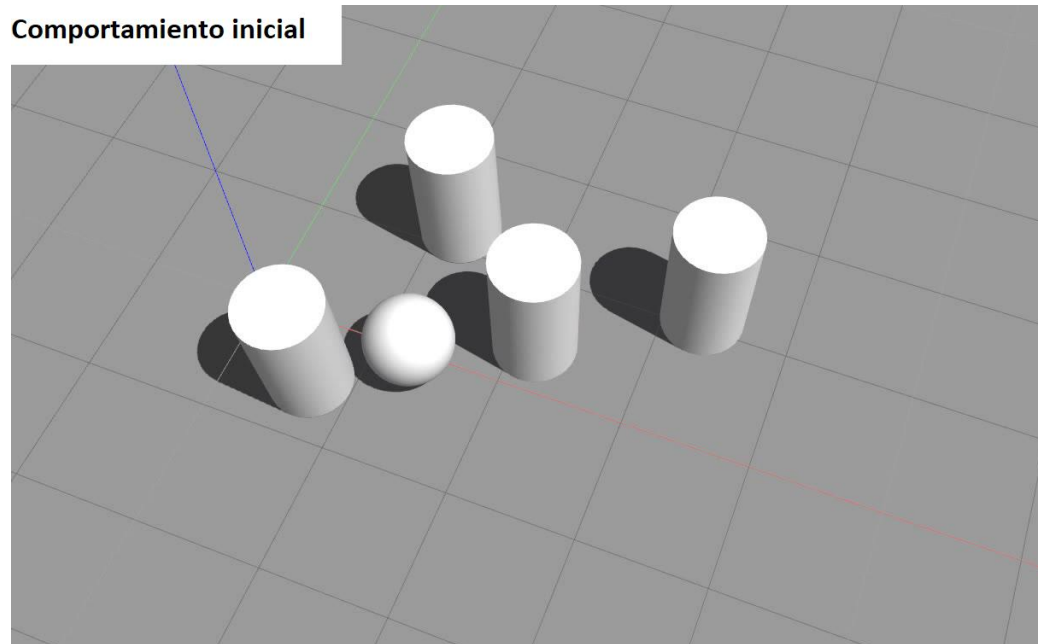


Generación mediante *splines*

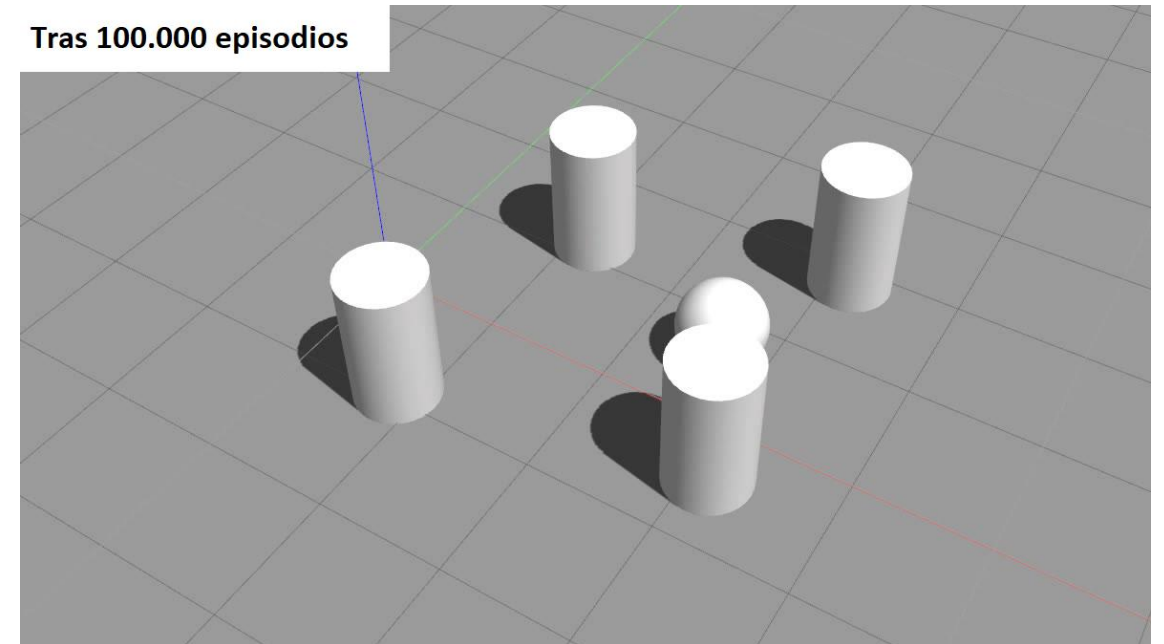


Planificador en simulación

COMPORTAMIENTO INICIAL



COMPORTAMIENTO TRAS 100000 EPISODIOS



Planificador en simulación

Sin *target networks*

Lazo abierto

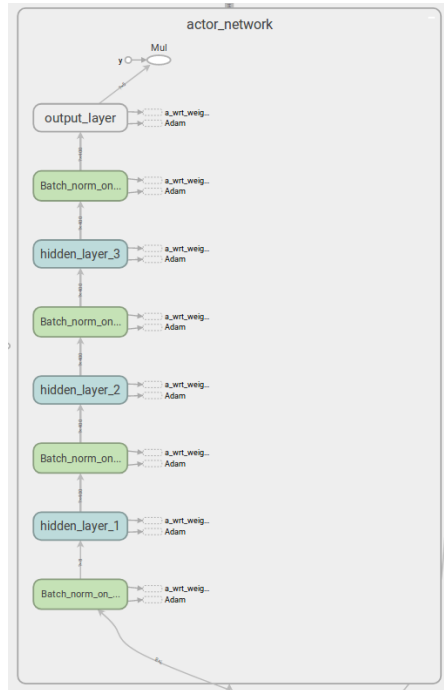
$$r_t = -o - 0,5d$$

Algoritmo 2: Algoritmo DDPG simplificado

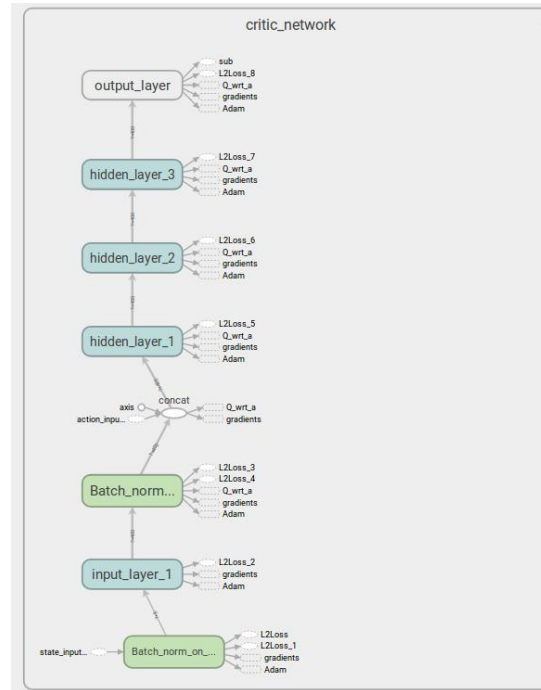
```
1 Inicializar aleatoriamente la red del crítico  $Q(s, a \mid \theta^Q)$  y el actor  
    $\mu(s \mid \theta^\mu)$  con pesos  $\theta^Q$  y  $\theta^\mu$ .  
2 Inicializar el buffer de interacciones  $R$   
3 for episodio=1,  $M$  do  
4   Inicializar proceso aleatorio  $N$  para la exploración del espacio de  
     acciones  
5   Recibir estado observacional inicial  $s_1$   
6   for  $t=1, T$  do  
7     Seleccionar acción  $a_t = \mu(s_t \mid \theta^\mu) + N_t$  de acorde con la política  
       actual y el ruido exploratorio  
8     Ejecutar acción  $a_t$  y observar la recompensa  $r_t$  y el nuevo  
       estado  $s_{t+1}$   
9     Almacenar transición  $(s_t, a_t, r_t, s_{t+1})$  en  $R$   
10    Muestrear un minibatch aleatorio de  $N$  transiciones  
       $(s_i, a_i, r_i, s_{i+1})$  de  $R$   
11    Sea  $y_i = r_i$   
12    Actualizar el crítico minimizando la pérdida:  
      
$$L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i \mid \theta^Q))^2$$
  
13    Actualizar la política del actor utilizando el gradiente de la  
      política muestreado:  
      
$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a \mid \theta^Q) \big|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s \mid \theta^\mu) \big|_{s_i}$$
  
14  end  
15 end
```

Planificador en simulación

ACTOR



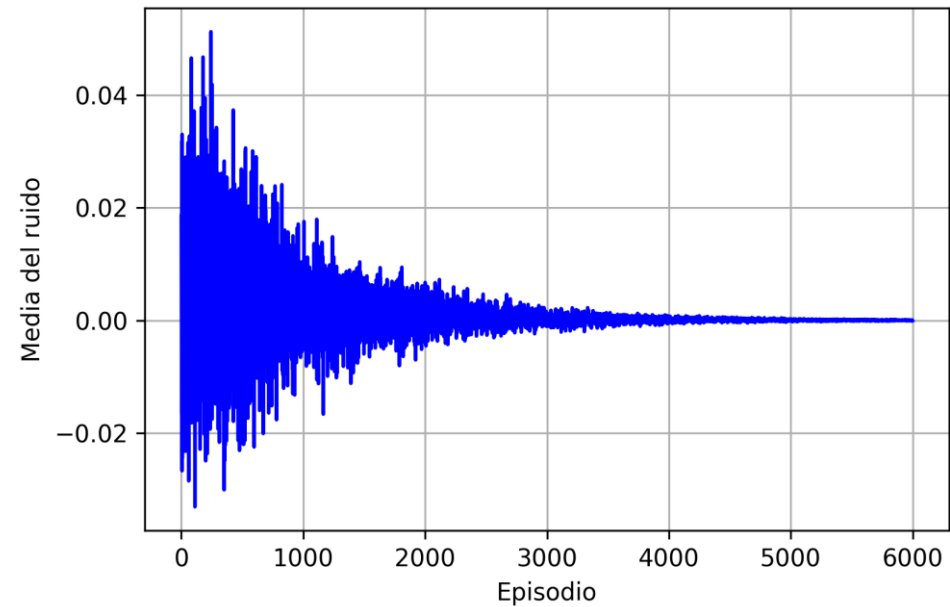
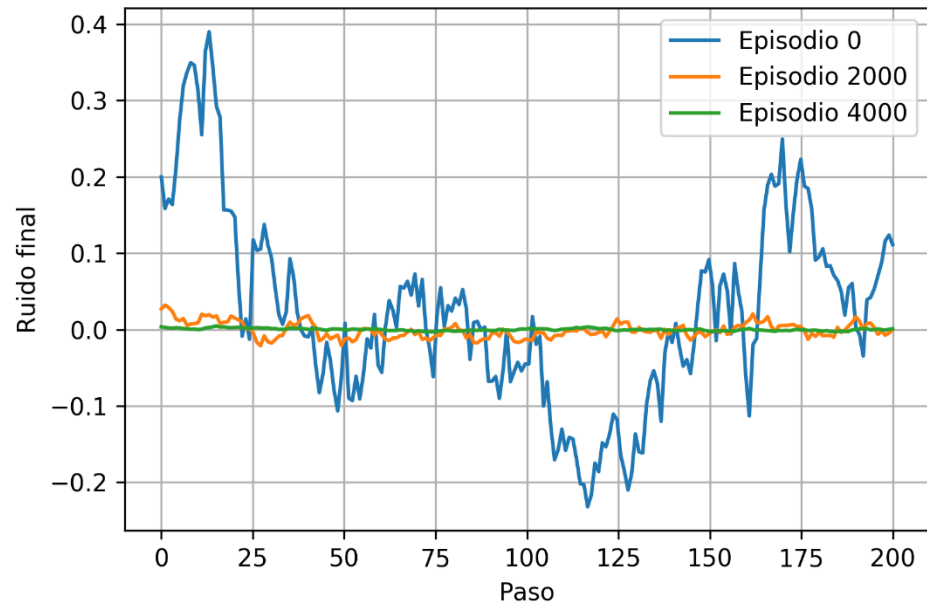
CRÍTICO



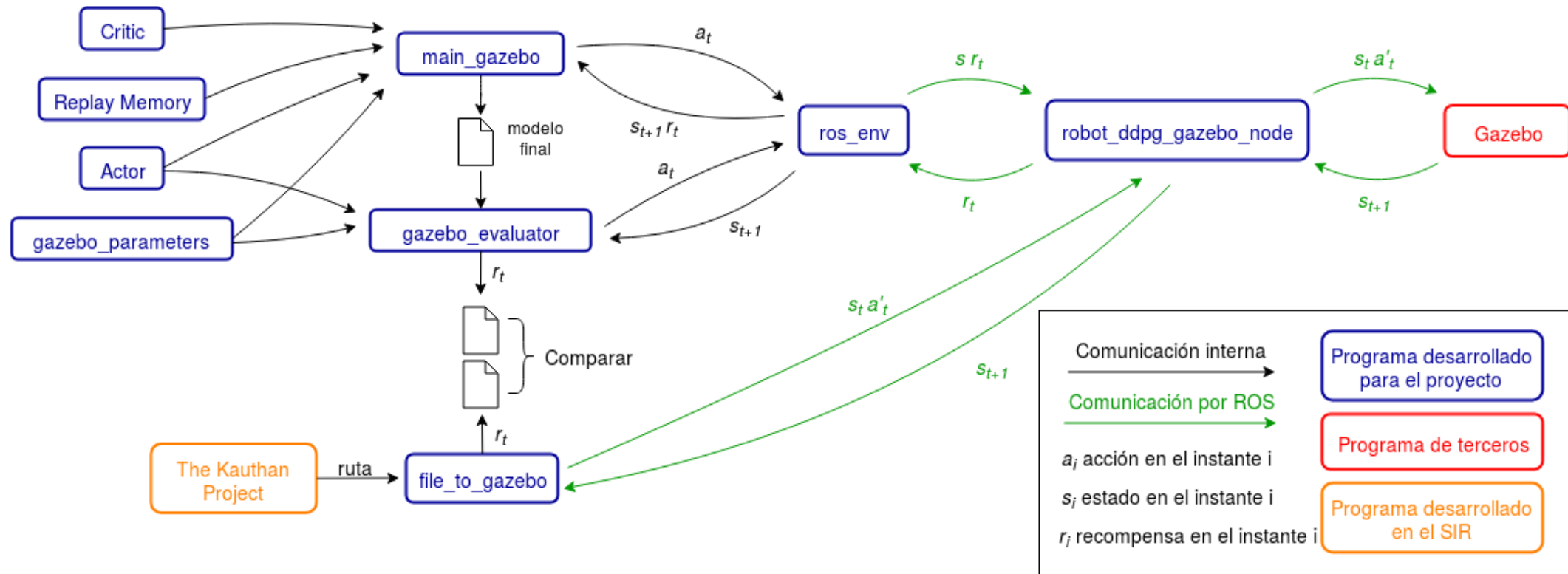
Elemento	Descripción
Arquitectura del crítico	2 capas internas de 400 unidades respectivamente. Todas las capas tienen activación ReLU.
Arquitectura del actor	2 capas internas de 400 unidades respectivamente. Todas las capas tienen activación ReLU excepto la salida, que usa una tanh.
Regularización L_2 weight decay para el crítico	10^{-2}
Optimizador para el crítico	ADAM
Optimizador para el actor	ADAM
Ratio de aprendizaje para el crítico	10^{-3}
Ratio de aprendizaje para el actor	10^{-4}
Factor de descuento γ	0.99
Factor τ para las actualizaciones de las <i>target networks</i>	0,01
Tamaño del <i>replay buffer</i>	10^5
Proceso Ornstein-Uhlenbeck	$\theta = 0,15, \sigma = 0,2$

Planificador en simulación

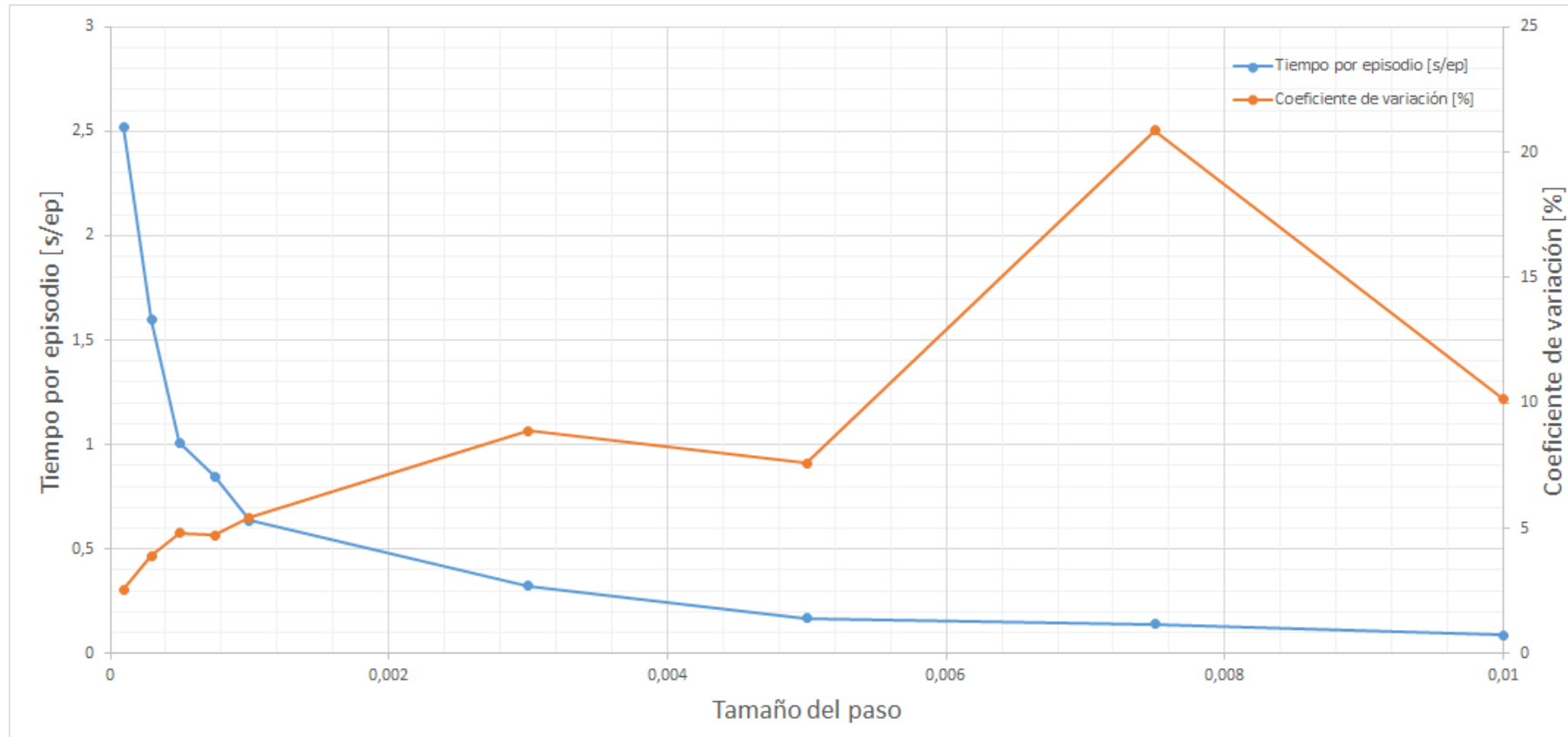
Proceso Ornstein-Uhlenbeck



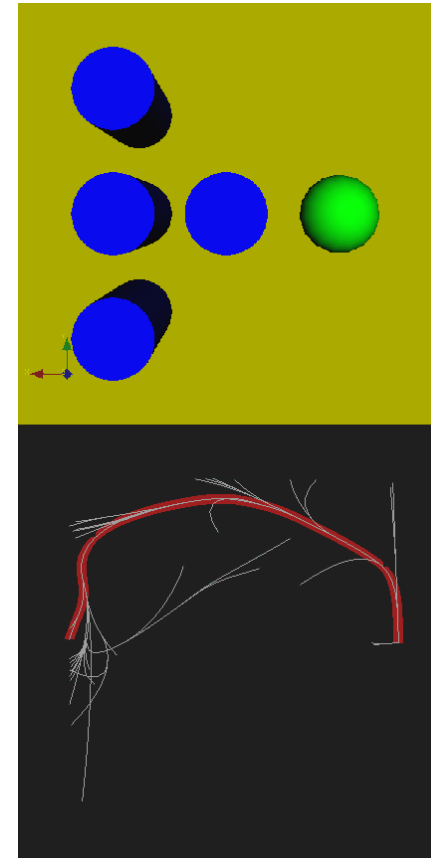
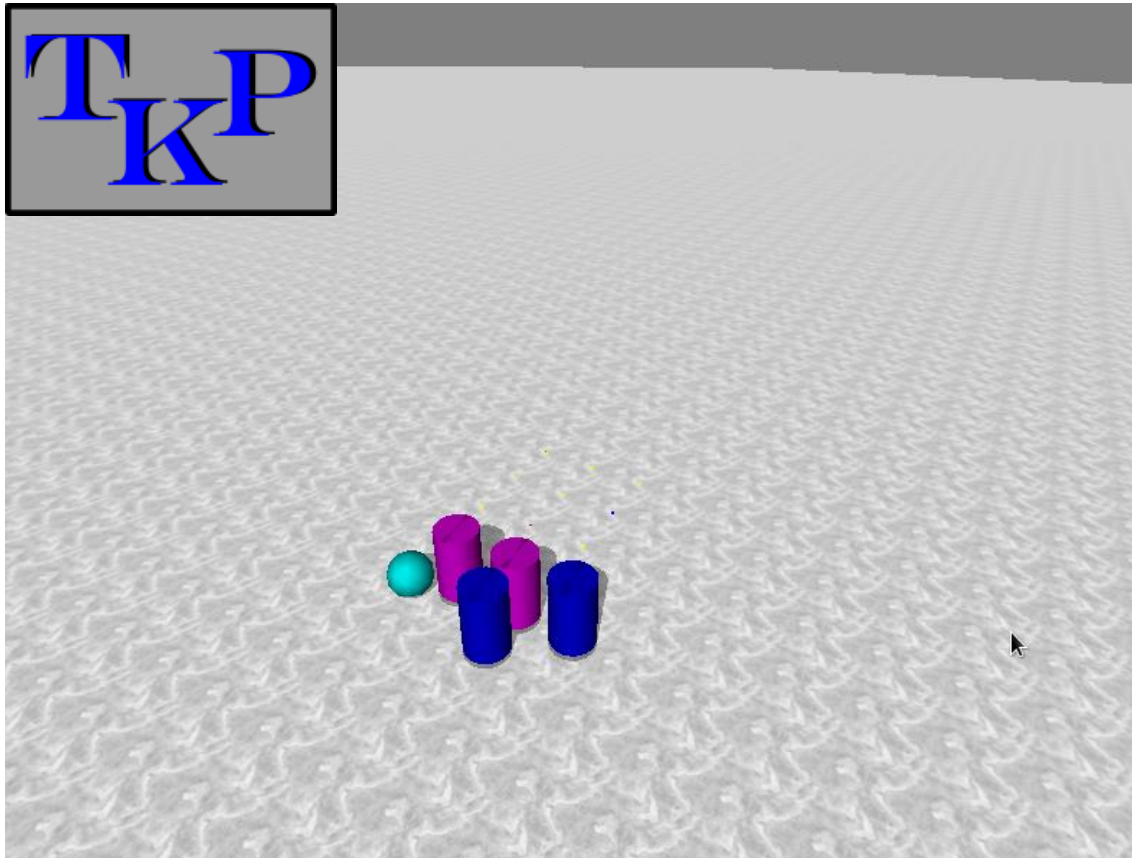
Planificador en simulación



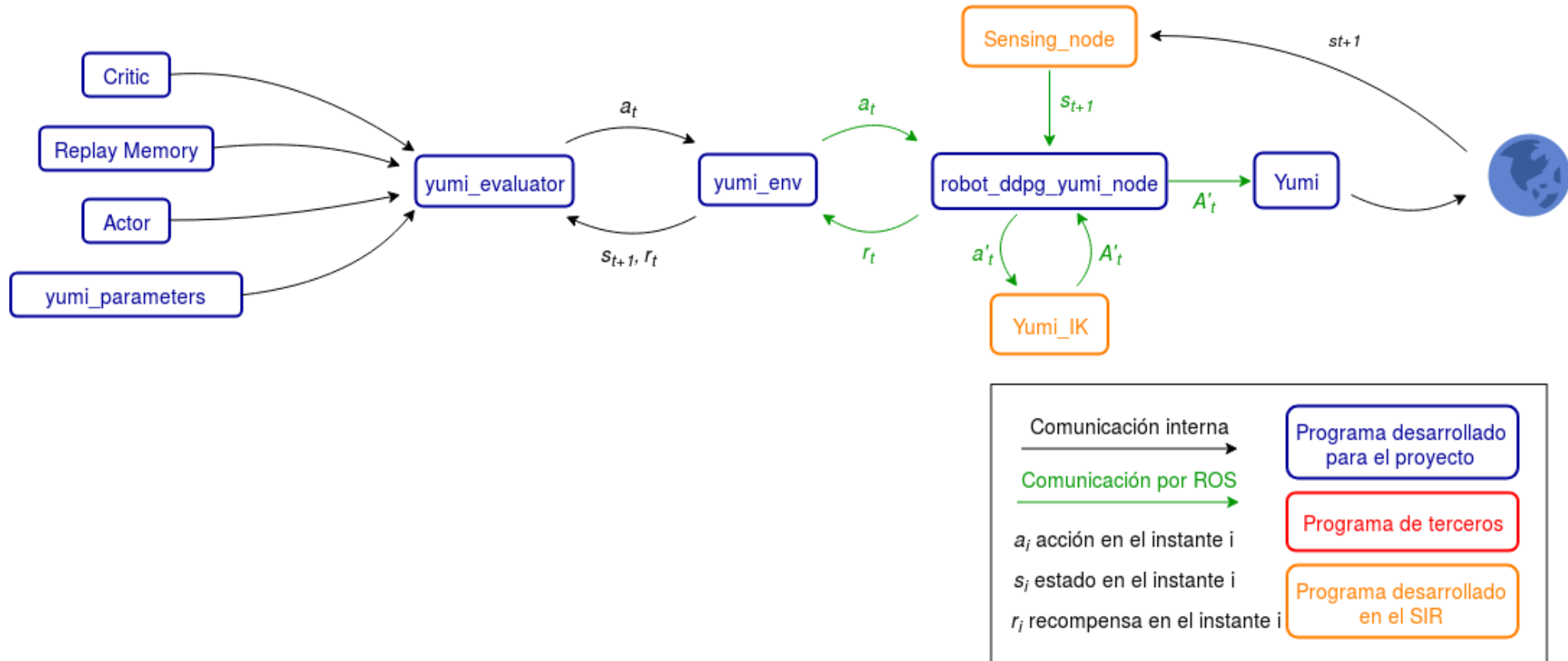
Planificador en simulación



The kautham Project: KPIECE



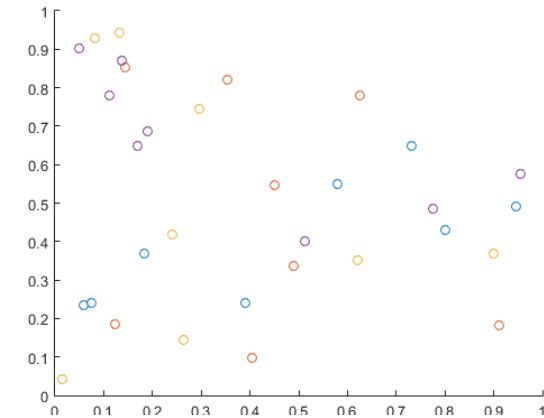
Implementación en YuMi



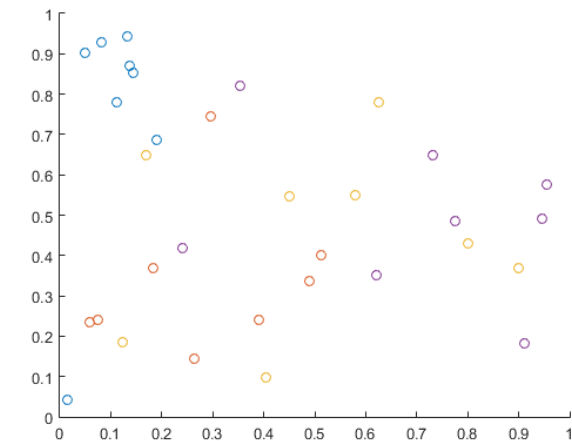
Experimentos

Experimento	Entornos	Descripción
Normal	Pendulum Mountain Car Simulación	DDPG básico, sin batch normalization.
Acciones en la segunda capa	Pendulum Mountain Car	Como en el algoritmo original, las acciones se añaden en la segunda capa del crítico.
<i>Batch normalization</i>	Pendulum Mountain Car	DDPG con batch normalization.
Sin <i>warmup</i>	Pendulum Mountain Car	Sin episodios de calentamiento.
Como el artículo original	Pendulum Mountain Car	DDPG tal y como se describe en el artículo original.
Entradas reordenadas	Simulación	DDPG básico, sin batch normalization pero con las entradas reordenadas.
Como el artículo original y con entradas reordenadas	Simulación	DDPG tal y como se describe en el artículo original y con las entradas reordenadas.

Original



Reordenado



Resultados

Introducción

- Objetivos
- Motivación

Algoritmo

- Antecedentes
- Puntos clave
- Núcleo

Implementación

- Entornos
- Arquitectura
- Detalles
- Experimentos

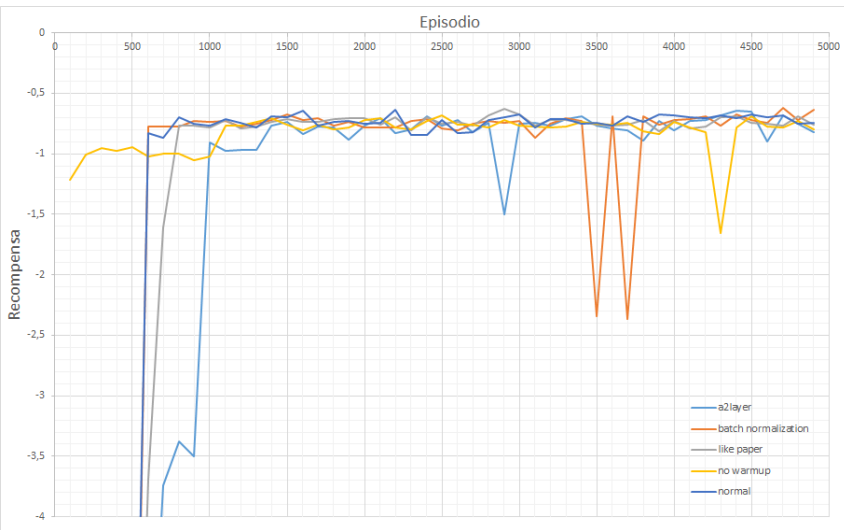
Resultados

- Pruebas con modificaciones
- Comparativa con KPIECE

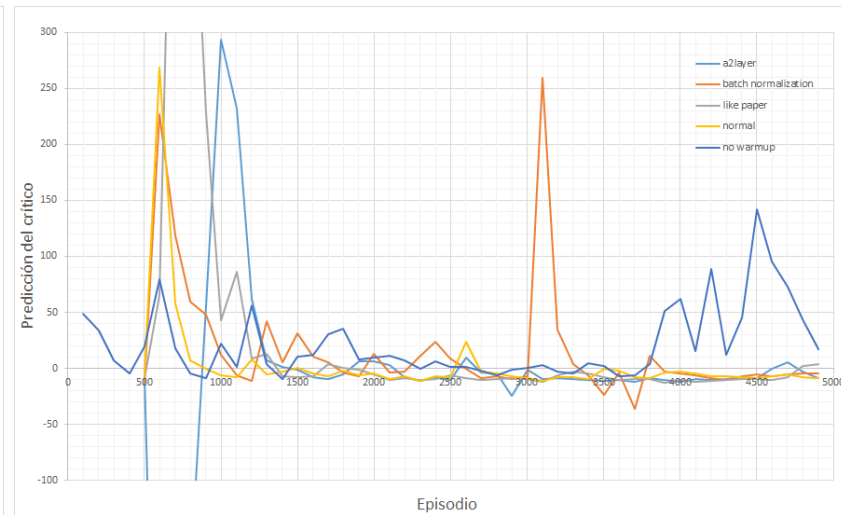
Conclusiones y trabajo futuro

Pendulum

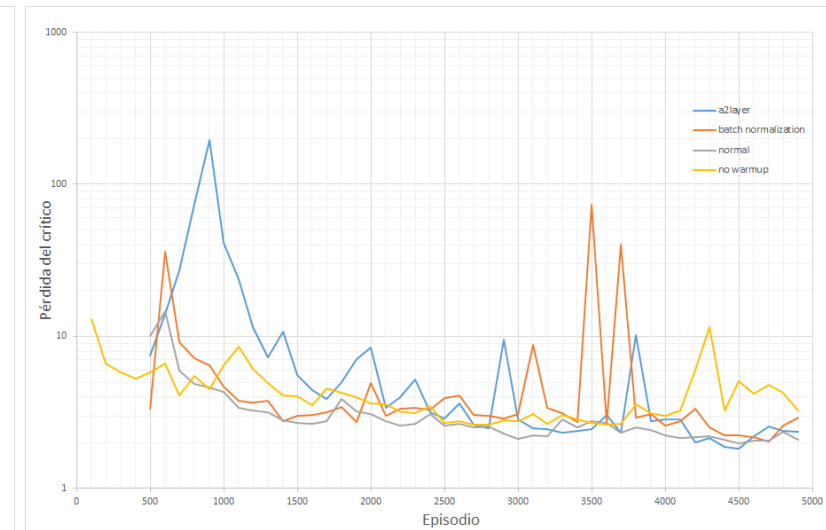
Recompensa



Predicción del crítico

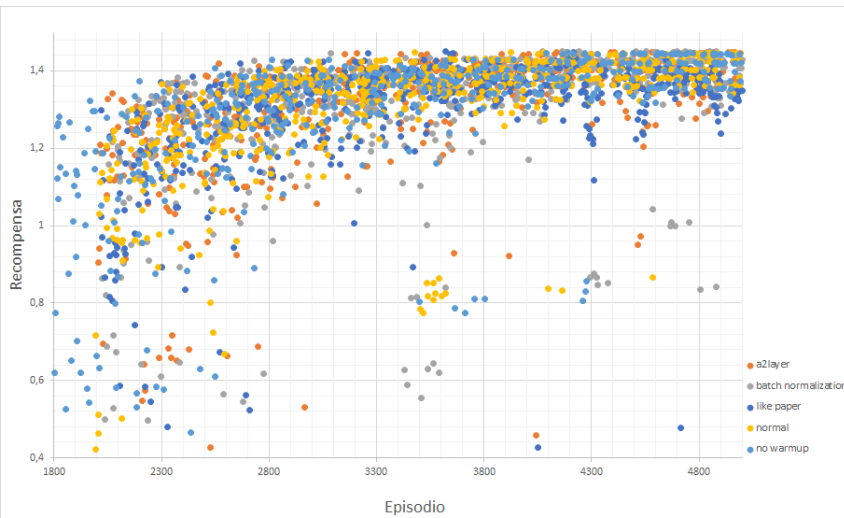


Pérdida del crítico

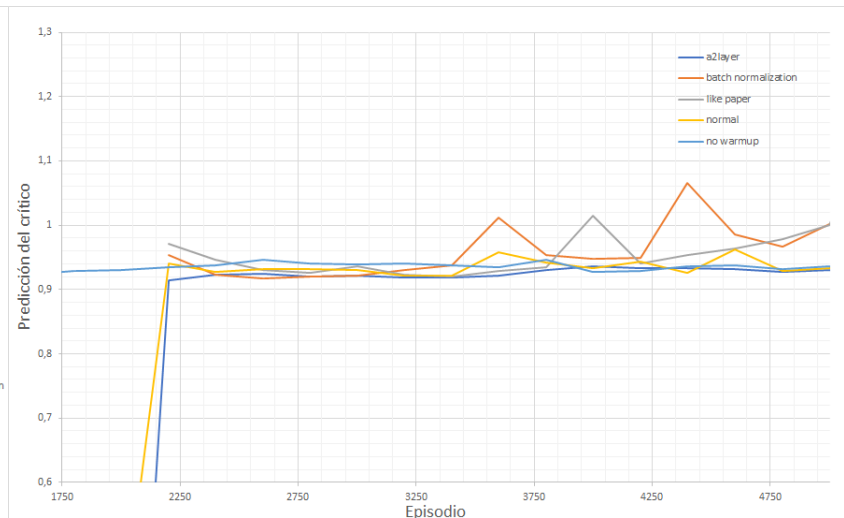


Mountain Car

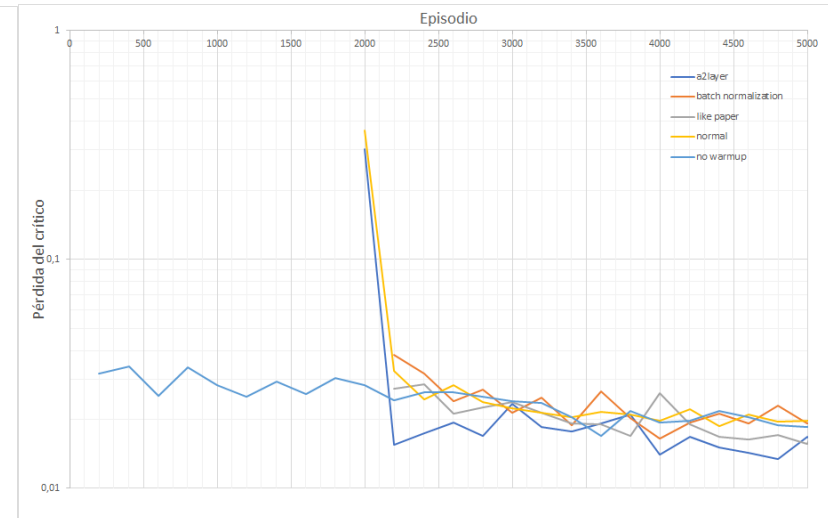
Recompensa



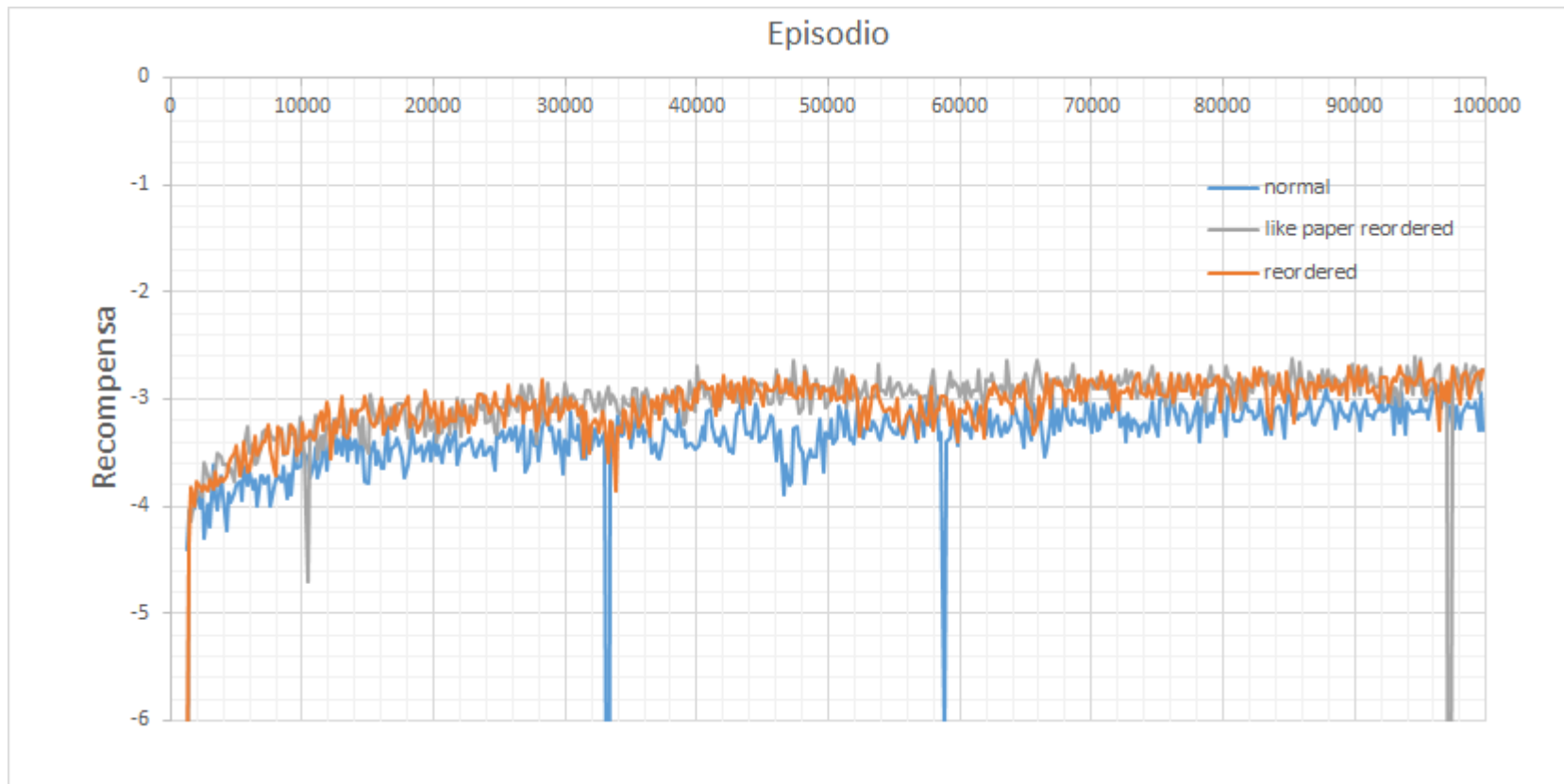
Predicción del crítico



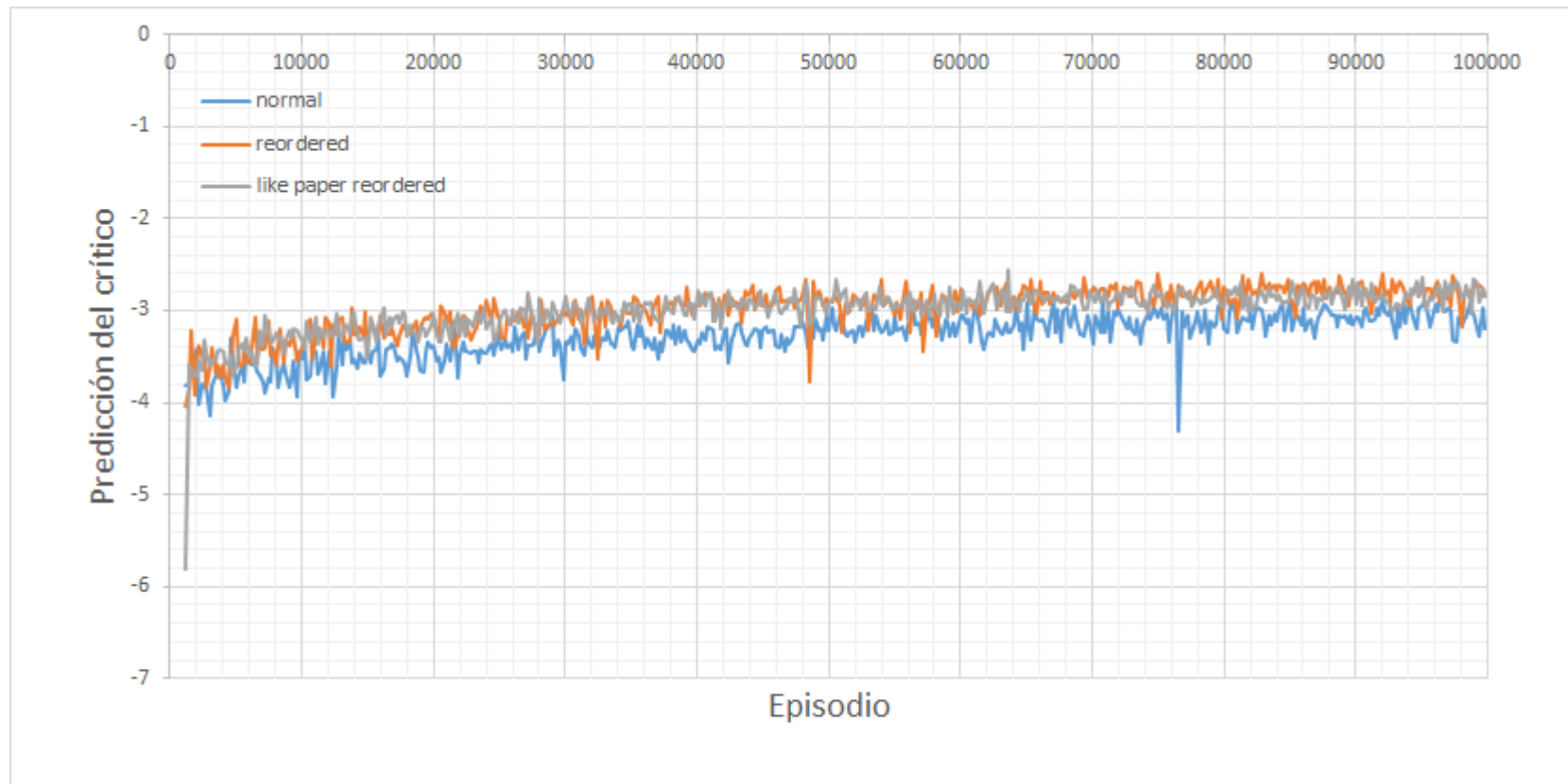
Pérdida del crítico



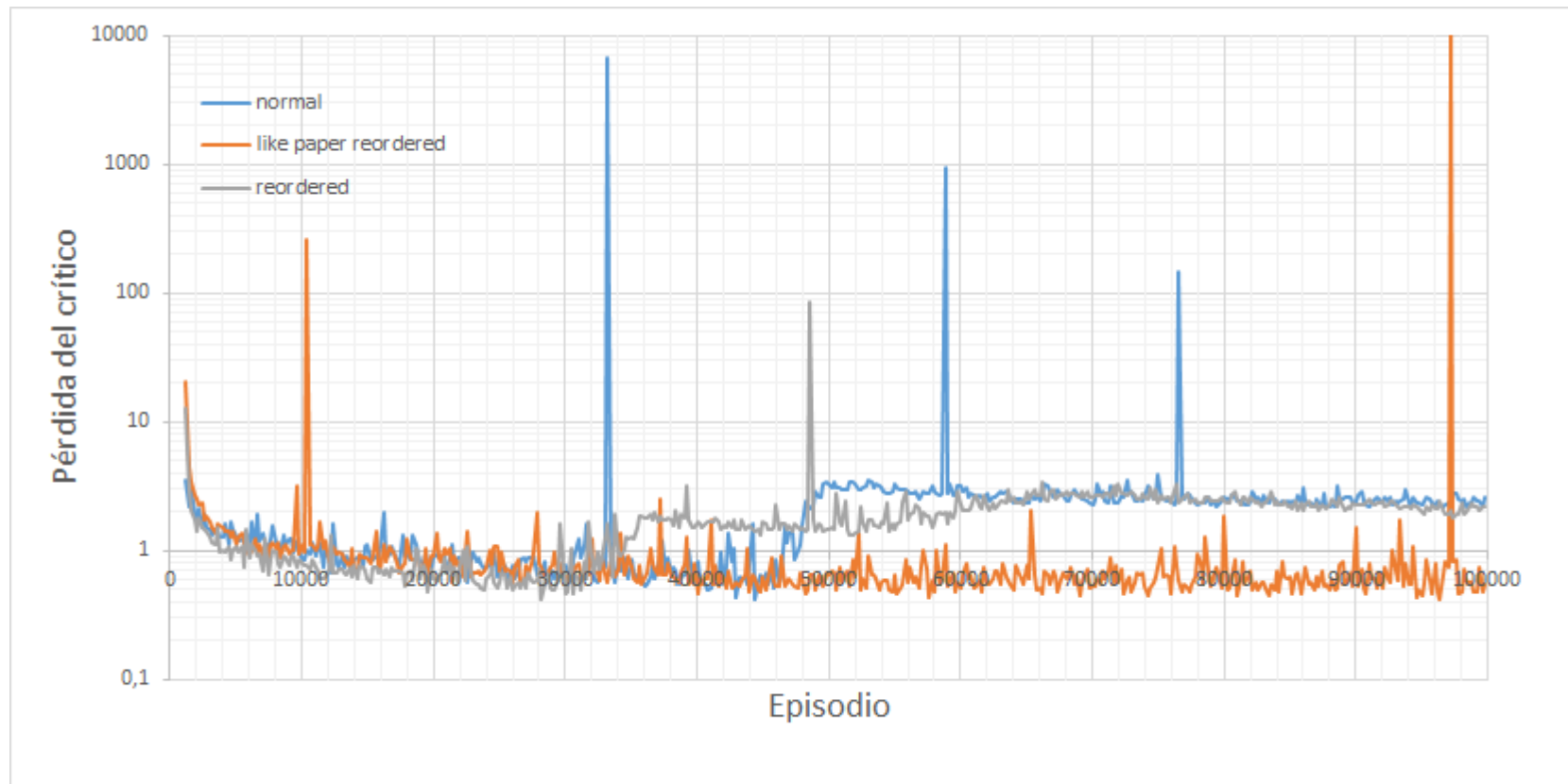
Planificador en simulación



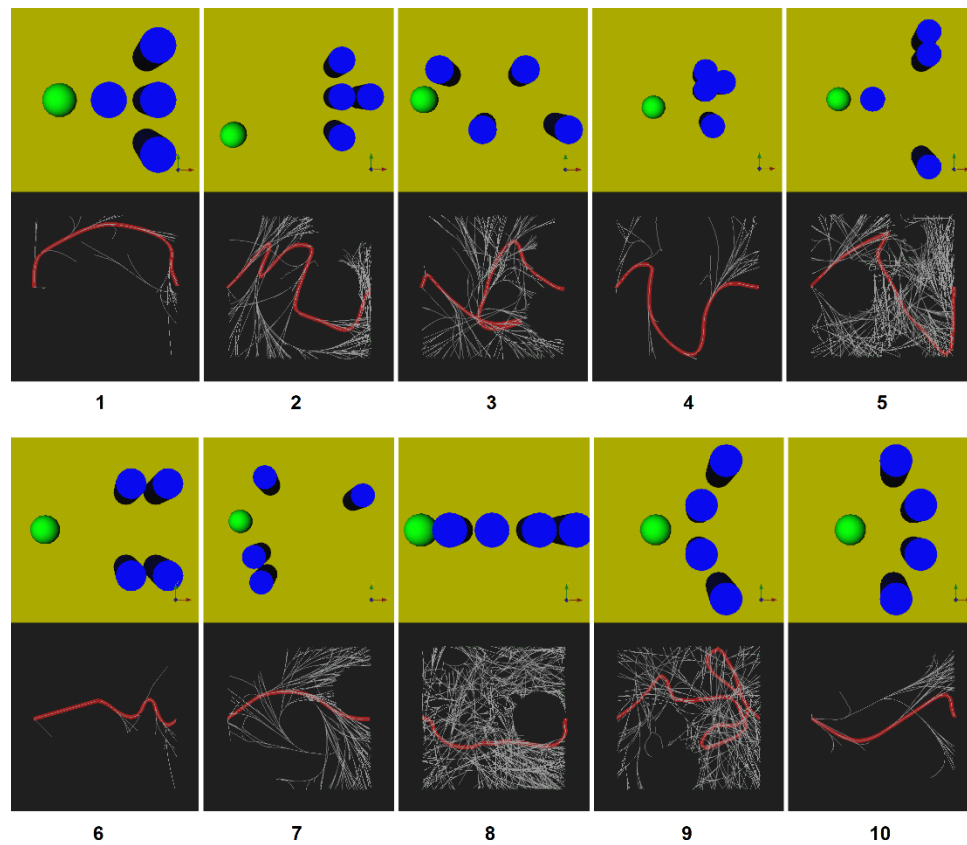
Planificador en simulación



Planificador en simulación

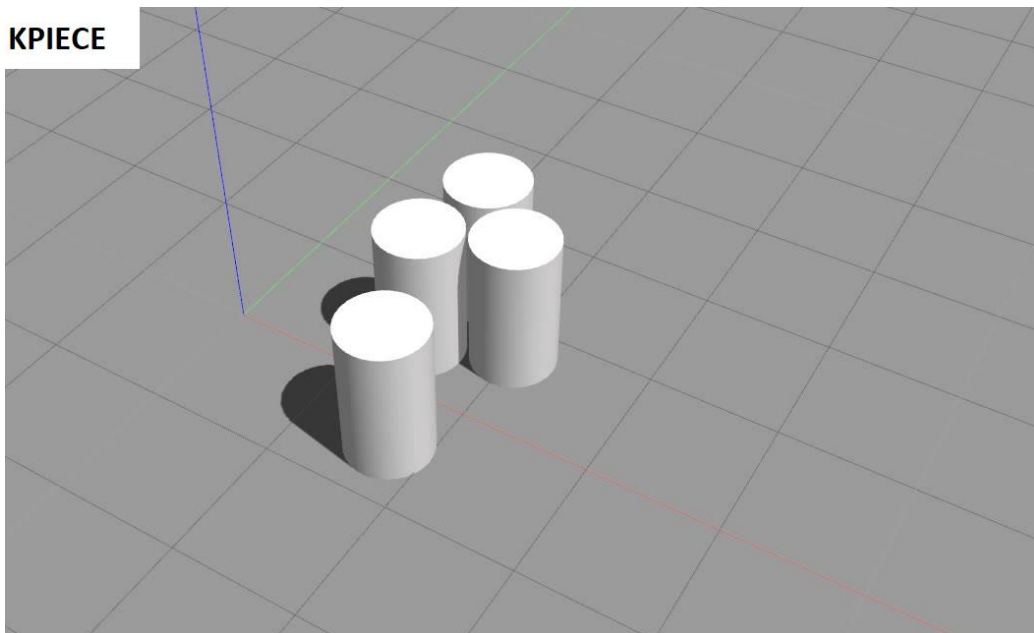


DDPG vs KPIECE

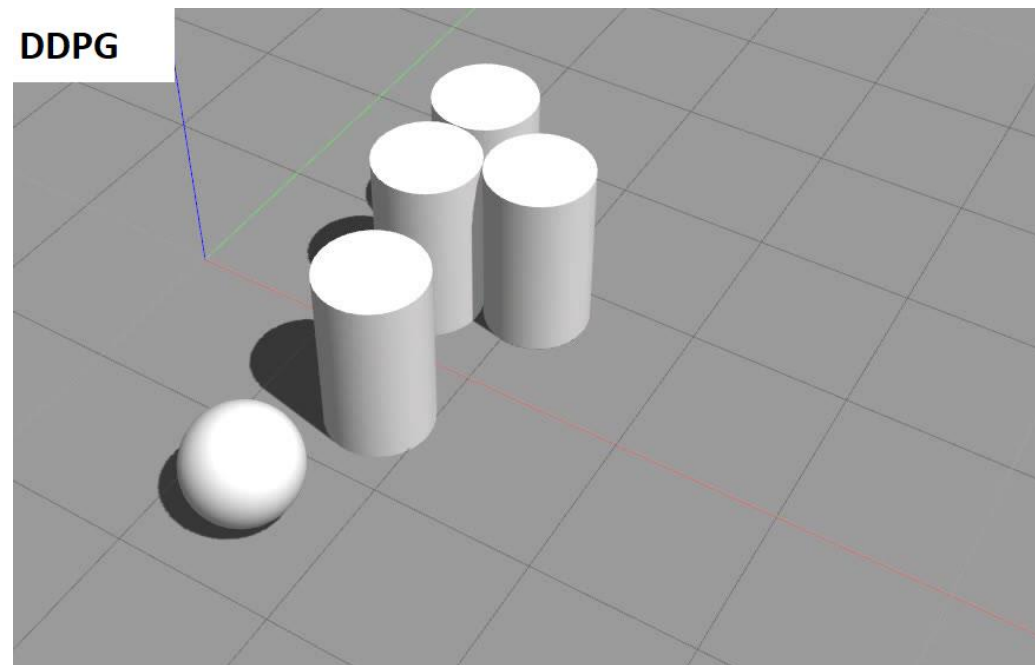


DDPG vs KPIECE

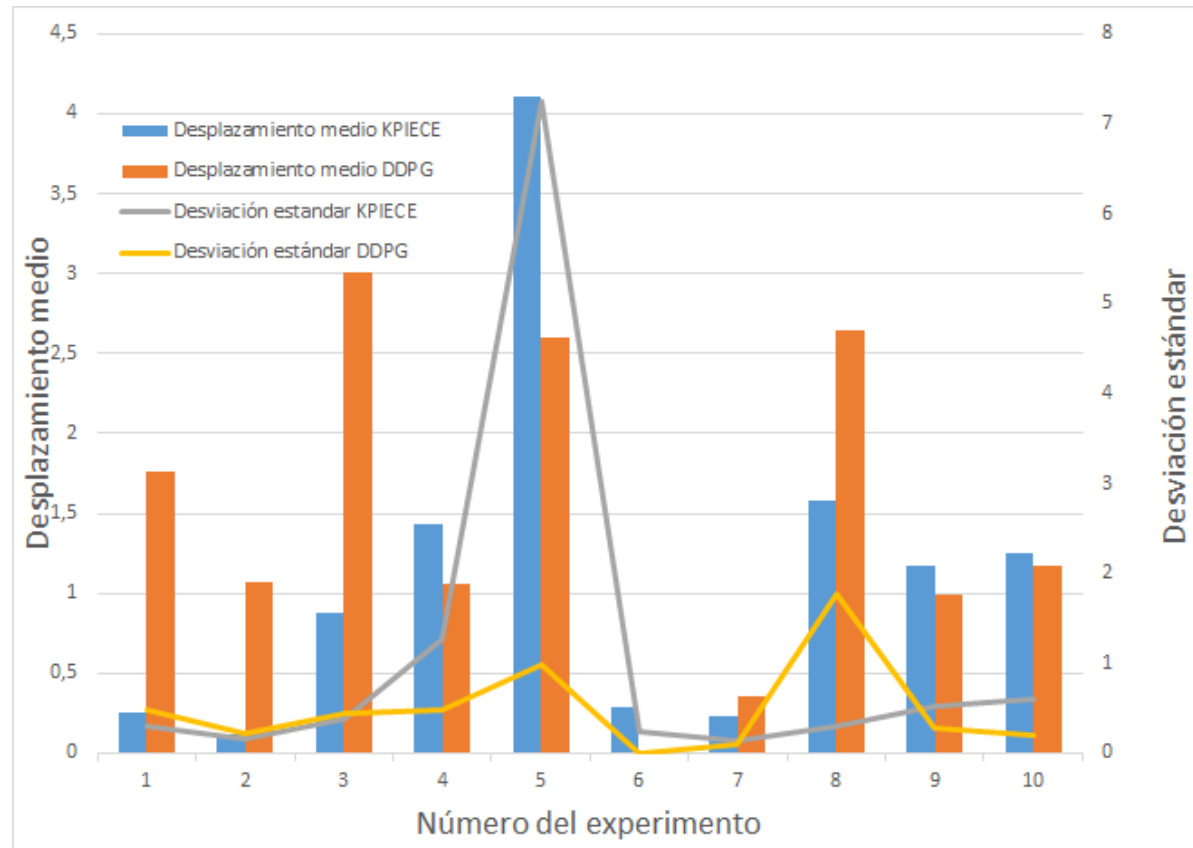
KPIECE



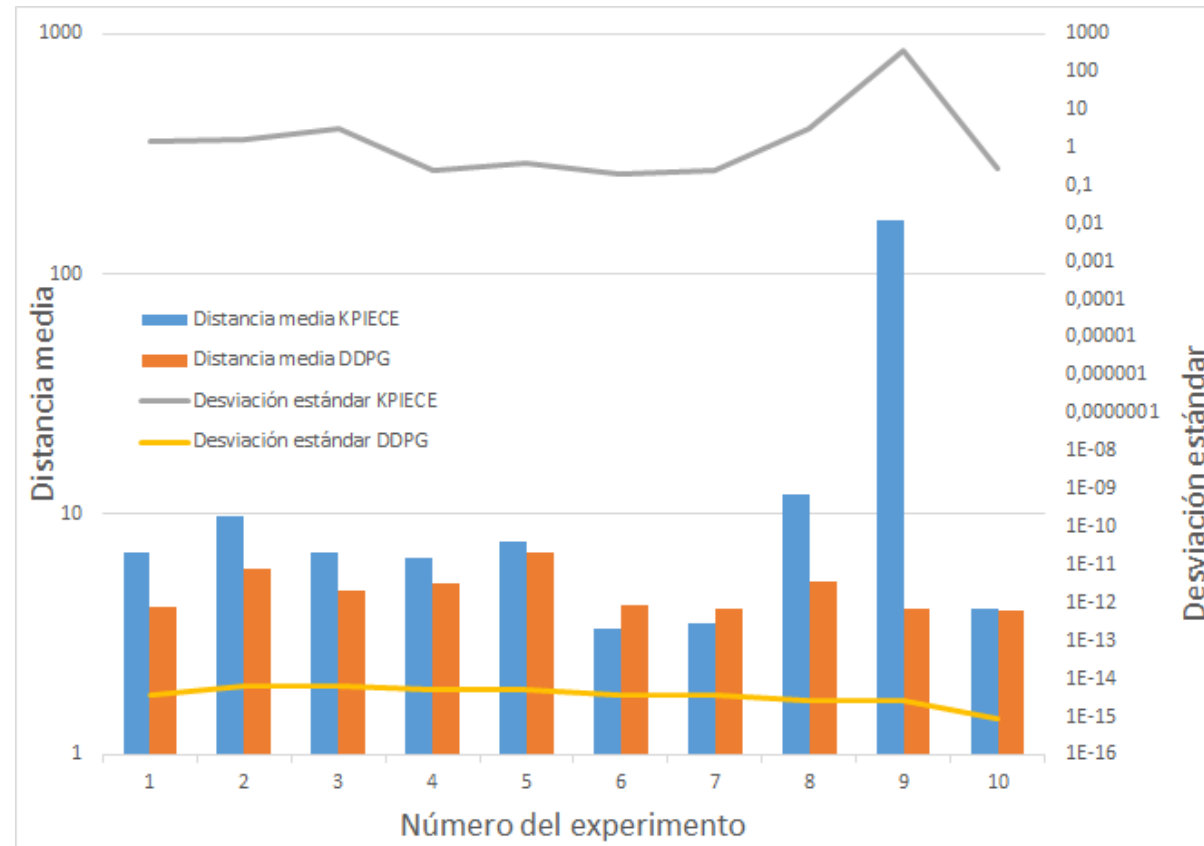
DDPG



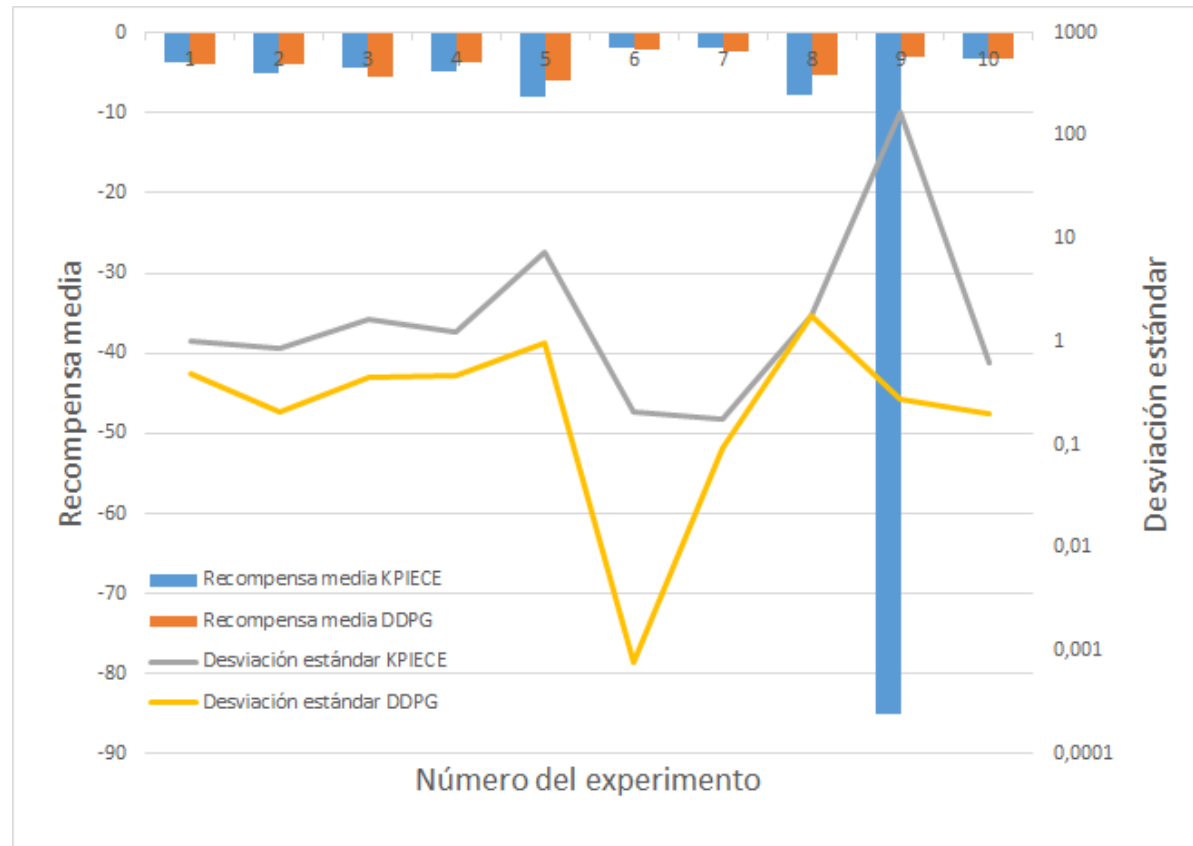
DDPG vs KPIECE



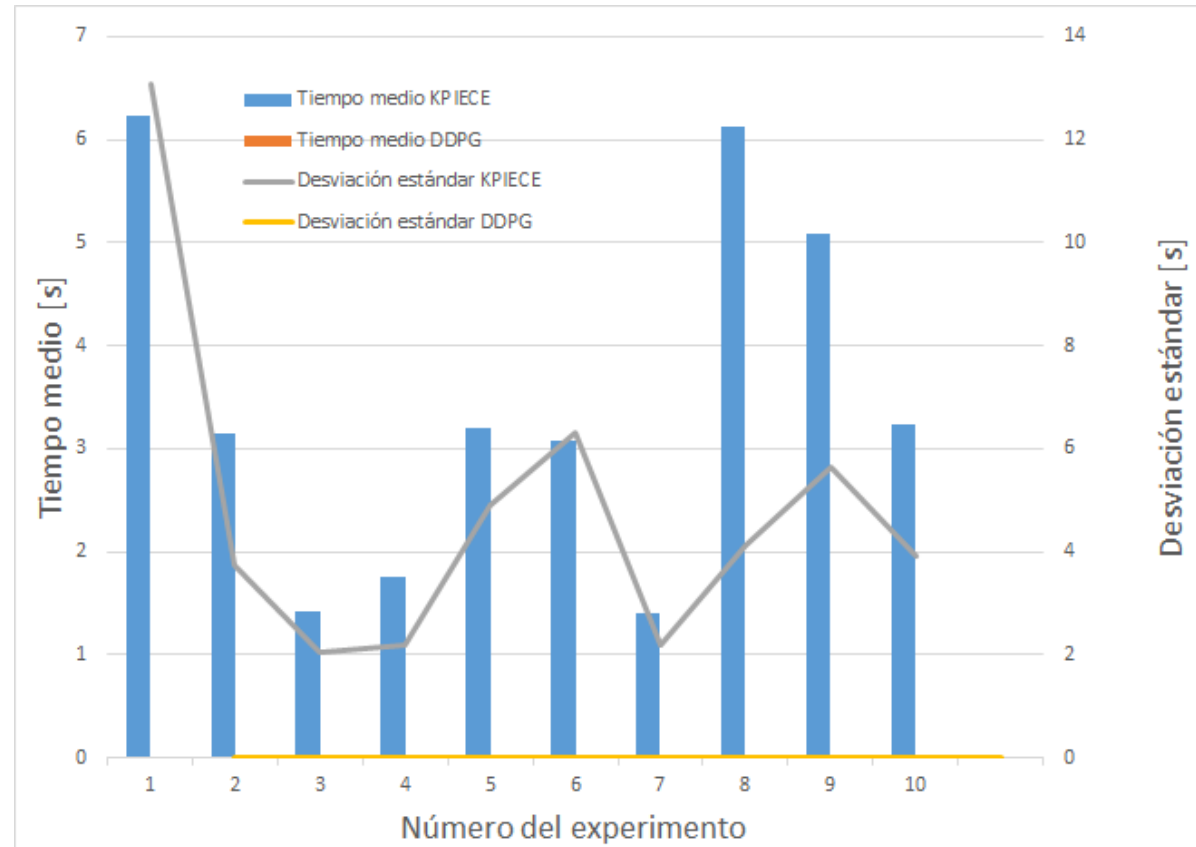
DDPG vs KPIECE



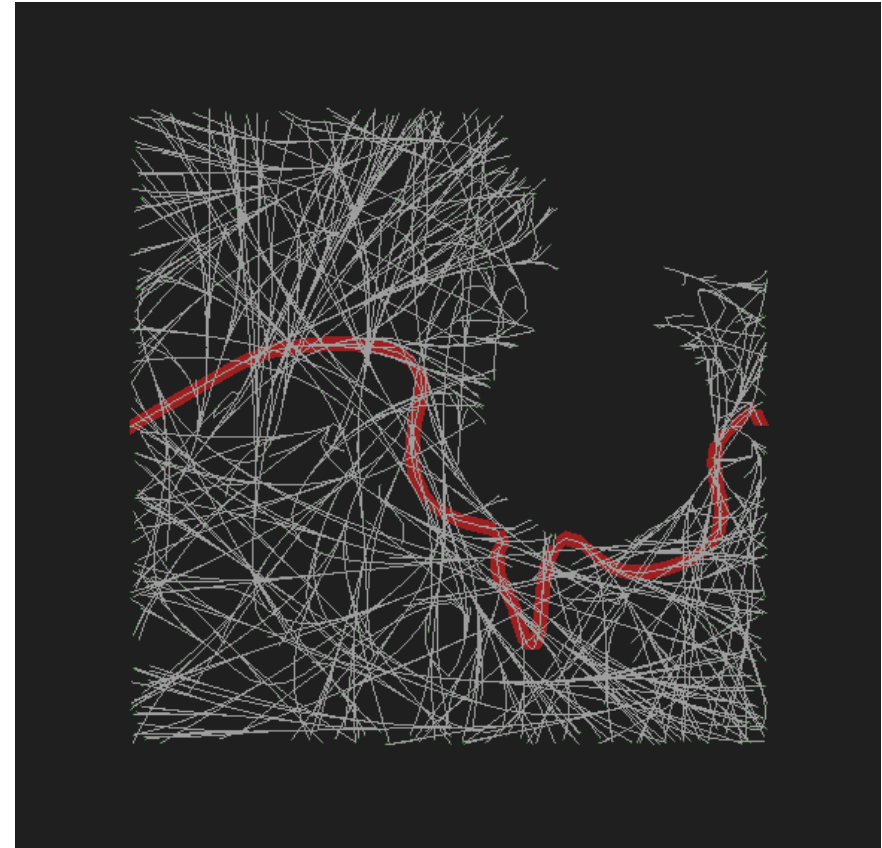
DDPG vs KPIECE



DDPG vs KPIECE



DDPG vs KPIECE



Conclusiones y trabajo futuro

Gran potencial

- Rapidez de respuesta
- Generalización
- Trabajo en entornos de alta dimensionalidad

Trabajo futuro

- Mejoras al *experience replay*
- Planificador en lazo cerrado

Muchas gracias por su atención